

1 We claim:

2 1. A computer-implemented method of text equivalencing from a string of characters

3 comprising:

4 modifying the string of characters using a predetermined set of heuristics;

5 comparing the modified string with a known string of characters in order to locate a

6 match;

7 responsive to not finding a match, forming a plurality of sub-strings of characters

8 from the string of characters; and

9 using an information retrieval technique on the sub-strings of characters to determine

8 a known string of characters equivalent to the string of characters.

2. The method of claim 1, wherein the information retrieval technique further comprises:

weighting the sub-strings;

scoring the known string of characters; and

retrieving information associated with the known string of characters with the highest

5 score.

1 3. The method of claim 2, further comprising, responsive to the highest score being

2 greater than a first threshold, automatically accepting the known string of characters as an exact

3 match.

1 4. The method of claim 2, further comprising, responsive to the highest score being less
2 than a second threshold and greater than a first threshold, presenting the known string of
3 characters to a user for manual confirmation.

1 5. The method of claim 2, further comprising, responsive to the highest score being less
2 than a second threshold and greater than a third threshold, presenting the known string of
3 characters to a user to select the equivalent string of characters.

6. The method of claim 1, wherein the sub-strings of characters are 3-grams.

7. The method of claim 1, wherein the string of characters is selected from the group
consisting of a song title, a song artist, an album name, a book title, an author's name, a book
publisher, a genetic sequence, and a computer program .

8. The method of claim 1, wherein the predetermined set of heuristics comprises
removing whitespace from the string of characters.

9. The method of claim 1, wherein the predetermined set of heuristics comprises
removing a portion of the string of characters.

10. The method of claim 1, wherein the predetermined set of heuristics comprises
replacing a symbol in the string of characters with an alternate representation for the symbol.

1 11. The method of claim 1 further comprising storing an indication that the string of
2 characters is the equivalent of the known string of characters.

1 12. A computer implemented system for text equivalencing from a string of characters
2 comprising:

3 a heuristics module for modifying the string of characters using a predetermined set
4 of heuristics;

5 a comparator module, coupled to the heuristics module, for comparing the modified
6 string with a known string of characters in order to find a match;

7 a sub-string formation module, coupled to the comparator module, responsive to not
8 finding a match, for forming a plurality of sub-strings of characters from the
9 string of characters; and

10 an information retrieval module, coupled to the sub-string formation module, for
11 performing an information retrieval technique on the sub-strings of characters
12 to determine a known string of characters equivalent to the string of
13 characters.

1 13. The system of claim 12, wherein the information retrieval module further comprises:
2 a weight module for weighting the sub-strings;
3 a score module for scoring the known string of characters; and

4 a retrieval module, coupled to the weight and score modules, for retrieving
5 information associated with the known string of characters with the highest
6 score.

1 14. The system of claim 13, further comprising an accept module, coupled to the
2 retrieval module, for accepting the information retrieved as an exact match for the highest score
3 greater than a first threshold.

1 15. The system of claim 13, further comprising an accept module, coupled to the
2 retrieval module, for presenting the information retrieved to a user for manual confirmation for
3 the highest score less than a first threshold and greater than a second threshold.

1 16. The system of claim 13, further comprising an accept module, coupled to the
2 retrieval module, for presenting the information retrieved to the user as a set of options for a user
3 to select for the highest score less than a second threshold and greater than a third threshold.

1 17. The system of claim 12, wherein the sub-strings of characters are 3-grams.

1 18. The system of claim 12, wherein the string of characters is selected from the group
2 consisting of a song title, a song artist, an album name, a book title, and author's name, a book
3 publisher, a genetic sequence, and a computer program.

1 19. The system of claim 12, wherein the predetermined set of heuristics comprises
2 removing whitespace from the string of characters.

1 20. The system of claim 12, wherein the heuristics module comprises a removal module
2 for removing a portion of the string of characters.

1 21. The system of claim 12, wherein the heuristics module comprises a replacement
2 module for replacing a symbol in the string of characters with an alternate representation for the
3 symbol.

1 22. The system of claim 12 further comprising a database update module for storing an
2 indication that the known string of characters is the equivalent of the known string of characters.

3 23. A computer-readable medium comprising computer-readable code for performing
4 text equivalencing from a string of characters comprising:

5 computer-readable code adapted to modify the string of characters using a
6 predetermined set of heuristics;

7 computer-readable code adapted to compare the modified string with a known string
8 of characters in order to locate a match;

9 computer-readable code, responsive to not finding a match, adapted to form a
10 plurality sub-strings of characters from the string of characters; and

11 computer-readable code adapted to use an information retrieval technique on the sub-
12 strings of characters to determine a known string of characters equivalent to
13 the string of characters.

1 24. The computer-readable medium of claim 23, wherein the information retrieval
2 technique further comprises:

3 computer-readable code adapted to weight the sub-strings;
4 computer-readable code adapted to score the known string of characters; and
5 computer-readable code adapted to retrieve information associated with the known
6 string of characters with the highest score.

1 25. The computer-readable medium of claim 24, further comprising computer-readable
2 code, responsive to the highest score being greater than a first threshold, adapted to automatically
3 accept the known string of characters as an exact match.

1 26. The computer-readable medium of claim 24, further comprising computer-readable
2 core, responsive the highest score being less than a second threshold and greater than a first
3 threshold, adapted to present the known string of characters to a user for manual confirmation.

1 27. The computer-readable medium of claim 24, further comprising computer-readable
2 code, responsive to the highest score being less than a second threshold and greater than a third
3 threshold, adapted to present the known string of characters to a user to select the equivalent
4 string of characters.

1 28. The computer-readable medium of claim 23, wherein the sub-strings of characters are
2 3-grams.

1 29. The computer-readable medium of claim 23, wherein the string of characters selected
2 from a group consisting of a song title, a song artist, an album name, a book title, an author's
3 name, a book publisher, a genetic sequence, and a computer program.

1 30. The computer-readable medium of claim 23, wherein the predetermined set of
2 heuristics comprises removing whitespace from the string of characters.

1 31. The computer-readable medium of claim 23, wherein the predetermined set of
2 heuristics comprises removing a portion of the string of characters.

1 32. The method of claim 23, wherein the predetermined set of heuristics comprises
2 replacing a symbol in the string of characters with an alternate representation for the symbol.

1 33. The computer-readable medium of claim 23 further comprising updating the known
2 string of characters to indicate the string of characters is the equivalent of the known string of
3 characters.

1 34. A computer-implemented system for performing text equivalencing from a string of
2 characters comprising:

3 a modifying means for modifying the string of characters using a predetermined set
4 of heuristics;

5 a comparator means for comparing the modified string with a known string of
6 characters in order to locate a match;

7 responsive to not finding a match, a formation means for forming a plurality sub-
8 strings of characters from the string of characters; and
9 an information retrieval means for determining a known string of characters
10 equivalent to the string of characters.

1 35. The system of claim 34, wherein the information retrieval means further comprises:
2 a weight means for weighting the sub-strings;
3 a score means for scoring the known string of characters; and
4 a retrieval means for retrieving information associated with the known string of
5 characters with the highest score.